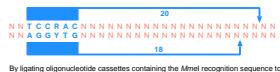


Restriction enzyme *Mme*I

Type IIa. A wonderful overhang cutter.



By fusing oligo-deoxycytidine containing the MmeI recognition sequence to DNA that has been cut with another restriction enzyme, called the anchoring enzyme (BamHI and MseI are just two possible choices), tags of length roughly 20 bp can be cut off, collected, and sequenced in large quantities.

20 bp allows 2^{20} , or 1.2×10^6 possible combinations. A single tag from a microbe to identify it to species level. Because it depends on specific locations relative to other genome features, even more information is added. Several tags can be used from an organism to identify species, in many cases.

Here are four schemes for locating tags. **Position reference features** are shown in red, anchor enzyme sites are shown in purple, tags are in blue.



3-most MseI site in mRNA transcript (SAGE, Velasco, p15, Proceedings from SAGE 2001: Proteins in Transcription Explosion, San Diego, 2001)



3-most MseI site in 16S rRNA (Bunn, et al., Genome Research 12, p1756, 2002)



MseI sites immediately upstream (3') of 16S rRNA standard primer (trials currently underway)



MseI sites immediately upstream (3') of heavily methylated regions (CpG islands).

Proof of principle: GSTs of a known genome

Genome Sequence Tags (GSTs) were collected from an avirulent strain of *Yersinia pestis* (EV76), for the purpose of comparing with the published genome sequence. The project is being done at Dunn, et al., Genome Research 12, p1756, 2002. GSTs were defined as 3-most MseI (CATG) and other sites in BamHI fragments.

Here's an example of a concatemer that was collected:

splices
The output from splices tag is a tabular matrix of tag sequences, described above and position. The common CATG profile in all tags is not shown. Tag lengths are evenly split between CATG +17 and +18; the variation is a product of the MmeI enzyme.

clust-tags
A simple histogram of tag frequency. The x-axis is tag sequence, the y-axis is frequency (count).

clust-tags
The output from clust-tags is a tabular matrix of tag sequences, described above and position. The common CATG profile in all tags is not shown. Tag lengths are evenly split between CATG +17 and +18; the variation is a product of the MmeI enzyme.

concatemers
Concatemer sequences were split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

1 Ambiguous tags are filtered out for rejection or separate processing

2 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

3 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

4 Ambiguous tags are filtered out for rejection or separate processing

5 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

6 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

7 Ambiguous tags are filtered out for rejection or separate processing

8 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

9 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

10 Ambiguous tags are filtered out for rejection or separate processing

11 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

12 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

13 Ambiguous tags are filtered out for rejection or separate processing

14 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

15 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

16 Ambiguous tags are filtered out for rejection or separate processing

17 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

18 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

19 Ambiguous tags are filtered out for rejection or separate processing

20 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

21 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

22 Ambiguous tags are filtered out for rejection or separate processing

23 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

24 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

25 Ambiguous tags are filtered out for rejection or separate processing

26 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

27 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

28 Ambiguous tags are filtered out for rejection or separate processing

29 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

30 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

31 Ambiguous tags are filtered out for rejection or separate processing

32 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

33 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

34 Ambiguous tags are filtered out for rejection or separate processing

35 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

36 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

37 Ambiguous tags are filtered out for rejection or separate processing

38 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

39 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

40 Ambiguous tags are filtered out for rejection or separate processing

41 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

42 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

43 Ambiguous tags are filtered out for rejection or separate processing

44 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

45 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

46 Ambiguous tags are filtered out for rejection or separate processing

47 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

48 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

49 Ambiguous tags are filtered out for rejection or separate processing

50 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

51 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

52 Ambiguous tags are filtered out for rejection or separate processing

53 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

54 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

55 Ambiguous tags are filtered out for rejection or separate processing

56 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

57 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

58 Ambiguous tags are filtered out for rejection or separate processing

59 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

60 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

61 Ambiguous tags are filtered out for rejection or separate processing

62 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

63 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAGAGGATGCGATATT
or
CATG + AACATGGGAGTGGCTATT

64 Ambiguous tags are filtered out for rejection or separate processing

65 Concatemer sequences are split into tags and processed with relatively simple perl scripts, shown in boxes. **splices** splits concatemer sequences into individual sequences, removes punctuation, and filters out tags which are too short or too long.

66 Three possible polarity punctuation conditions must be considered:

GTGGGGCGGCCGCTATT
Forward tag: CATG + GTGGGGCGGCCGCTATT

AATGTTGATGAGATGTC
Reverse tag: CATG + GGGATGTTGATTACCA

AACAGAGGATGCGATATT
Ambiguous (Tag begins with AA in genome)
could be
CATG + AACAG